



Accompanying clean version (without markings) of the Substitute Specification

PATENT

10830.0098.NPUS00

APPLICATION FOR UNITED STATES LETTERS PATENT

for

**MANAGEMENT OF THE FILE-MODIFICATION TIME ATTRIBUTE IN
A MULTI-PROCESSOR FILE SERVER SYSTEM**

By

Stephen A. Fridella

Gang Ma

Xiaoye Jiang

Sorin Faibish

Rui Liang

Express Mail Mailing Label No. EV 318623741 US

BACKGROUND OF THE INVENTION

1. Field of the Invention

The present invention relates generally to data storage systems, and more particularly to network file servers. The present invention specifically relates to a file server system in which access to file attributes is shared among a number of processors.

2. Description of the Related Art

Network data storage is most economically provided by an array of low-cost disk drives integrated with a large semiconductor cache memory. A number of data mover computers are used to interface the cached disk array to the network. The data mover computers perform file locking management and mapping of the network files to logical block addresses of storage in the cached disk array, and move data between network clients and the storage in the cached disk array.

Data consistency problems may arise if concurrent client access to a read/write file is permitted through more than one data mover. These data consistency problems can be solved in a number of ways. For example, as described in Vahalia et al., U.S. Patent 5,893,140 issued April 6, 1999, entitled "File Server Having a File System Cache and Protocol for Truly Safe Asynchronous Writes," incorporated herein by reference, locking information can be stored in the cached disk array, or cached in the data mover computers if a cache coherency scheme is used to maintain consistent locking data in the caches of the data mover computers.

When a large number of clients are concurrently accessing shared read-write files, there may be considerable access delays due to contention for locks not only on the files but also on the file directories. One way of reducing this contention is to assign each file

1 system to only one data mover assigned the task of managing the locks on the files and
2 directories in the file system. This permits the data mover file manager to locally cache
3 and manage the metadata for the files and directories of the file system. For example, as
4 described in Xu et al., U.S. Patent 6,324,581, issued Nov. 27, 2001, incorporated herein
5 by reference, the data mover acting as the manager of a file grants a lock on the file and
6 provides metadata of the file to another data mover servicing a client request for access to
7 the file. Then the data mover servicing the client request uses the metadata to directly
8 access the file data in the cached disk array.

9 It is desired to permit clients to have asynchronous writes to a file in accordance
10 with version 3 of the Network File System (NFS) protocol, and concurrent write access
11 and byte range locking to a file in accordance with version 4 of the NFS protocol. (See
12 NFS Version 3 Protocol Specification, RFC 1813, Sun Microsystems, Inc., June 1995,
13 incorporated herein by reference, and NFS Version 4 Protocol Specification, RFC 3530,
14 Sun Microsystems, Inc., April 2003, incorporated herein by reference.) In this case, it is
15 possible for a file to be updated at about the same time by multiple clients. The NFS
16 protocol specifies that the time of last update of a file should be indicated by a file-
17 modification time attribute, referred to in the protocol as "mtime."

18 19 **SUMMARY OF THE INVENTION**

20 In accordance with one aspect, the invention provides a method of operation in a
21 file server system. The file server system has a clock for producing a clock time and a
22 processor for servicing client requests for access to a file. The processor has a timer for
23 measuring a time interval. The method includes the processor obtaining the clock time

1 from the clock, and beginning measurement of the time interval with the timer. The
2 method further includes the processor responding to a request from a client for an
3 asynchronous write to the file by performing an asynchronous write operation with
4 respect to the file, and determining a file-modification time that is a function of the clock
5 time having been obtained from the clock and the time interval measured by the timer,
6 the file-modification time indicating a time of modification of the file by the
7 asynchronous write operation.

8 In accordance with another aspect, the invention provides a method of operation
9 in a file server system having a first processor and a second processor for servicing client
10 requests for access to a file. The first processor has a clock producing a clock time, and
11 the second processor has a timer for measuring a time interval. The method includes the
12 second processor responding to a first request from a client for an asynchronous write to
13 the file by obtaining the clock time from the clock of the first processor, beginning
14 measurement of the time interval with the timer, performing a first asynchronous write
15 operation with respect to the file, and using the clock time obtained from the clock of the
16 first processor as a first file-modification time indicating a time of modification of the file
17 by the first asynchronous write operation. The method further includes the secondary
18 processor responding to a second request from the client for an asynchronous write to the
19 file by performing a second asynchronous write operation with respect to the file, and
20 determining a second file-modification time that is a function of the clock time obtained
21 from the clock of the first processor and the time interval measured by the timer. The
22 second file-modification time indicates a time of modification of the file by the second
23 asynchronous write operation.

1 In accordance with yet another aspect, the invention provides a method of
2 operation in a file server system having a first processor and a second processor for
3 servicing client requests for access to a file. The first processor has a clock producing a
4 clock time, and the second processor has a timer for measuring a time interval. The
5 method includes the second processor responding to a first request from a client for an
6 asynchronous write to the file by obtaining the clock time from the clock of the first
7 processor, beginning measurement of the time interval with the timer, performing a first
8 asynchronous write operation with respect to the file, and using the clock time obtained
9 from the clock of the first processor as a first file-modification time indicating a time of
10 modification of the file by the first asynchronous write operation. The method further
11 includes the second processor receiving from the first processor an updated value for the
12 file-modification time, the second processor comparing the updated value for the file-
13 modification time to the first file-modification time, and upon finding that the updated
14 value for the file-modification time is greater than the first file-modification time, the
15 second processor resetting the timer. Moreover, the method further includes the second
16 processor responding to a second request from the client for an asynchronous write to the
17 file by performing a second asynchronous write operation with respect to the file, and
18 determining a second file-modification time that is a function of the updated value for the
19 file-modification time and the time interval measured by the timer. The second file-
20 modification time indicates a time of modification of the file by the second asynchronous
21 write operation.

22 In accordance with yet another aspect, the invention provides a method of
23 operation in a file server system having a primary processor managing metadata of a file,

1 and a secondary processor responding to requests from a client for access to the file. The
2 primary processor has a clock producing a clock time, and the secondary processor has a
3 timer for measuring a time interval. The method includes the secondary processor
4 responding to a first asynchronous write request from the client for writing to the file by
5 obtaining attributes of the file and the clock time from the primary processor, storing the
6 attributes of the file in a cache local to the secondary processor and using the file
7 attributes to perform a first asynchronous write operation with respect to the file,
8 beginning measurement of the time interval with the timer, and using the clock time as a
9 first file-modification time indicating a time of modification of the file by the first
10 asynchronous write operation. The method further includes the secondary processor
11 responding to a second asynchronous write request from the client for writing to the file
12 by using the attributes of the file in the cache local to the secondary processor to perform
13 a second asynchronous write operation with respect to the file, and determining a second
14 file-modification time that is a function of the clock time having been obtained from the
15 clock of the primary processor and the time interval measured by the timer, the second
16 file-modification time indicating a time of modification of the file by the second
17 asynchronous write operation.

18 In accordance with still another aspect, the invention provides a method of
19 operation in a network file server. The network file server has a plurality of data mover
20 computers for servicing client requests for access to a file, and a cached disk array for
21 storing data of the file. The data mover computers are coupled to the cache disk array for
22 accessing the data of the file. The data mover computers include a primary data mover
23 computer managing metadata of the file, and a secondary data mover computer that

1 requests metadata of the file from the primary data mover computer. The primary data
2 mover computer has a clock producing a clock time, and the secondary data mover
3 computer has a timer for measuring a time interval. The method includes the secondary
4 data mover computer responding to a first asynchronous write request from a client for
5 writing to the file by obtaining attributes of the file and the clock time from the primary
6 data mover computer, storing the attributes of the file in a cache local to the secondary
7 data mover computer and using the file attributes to perform a first asynchronous write
8 operation with respect to the file, beginning measurement of the time interval with the
9 timer, and using the clock time as a first file-modification time indicating a time of
10 modification of the file by the first asynchronous write operation. The method further
11 includes the secondary data mover computer responding to a second asynchronous write
12 request from the client for writing to the file by using the attributes of the file in the cache
13 local to the secondary data mover computer to perform a second asynchronous write
14 operation with respect to the file, and determining a second file-modification time as a
15 function of the clock time having been obtained from the primary data mover and the
16 time interval measured by the timer, the second file-modification time indicating a time
17 of modification of the file by the second asynchronous write operation.

18 In accordance with another aspect, the invention provides a file server system
19 having a clock for producing a clock time and a processor for servicing client requests for
20 access to a file. The processor has a timer for measuring a time interval. The processor
21 is programmed for obtaining the clock time from the clock, and beginning measurement
22 of the time interval with the timer. The processor is further programmed for responding
23 to a request from a client for an asynchronous write to the file by performing an

1 asynchronous write operation with respect to the file, and determining a file-modification
2 time that is a function of the clock time having been obtained from the clock and the time
3 interval measured by the timer, the file-modification time indicating a time of
4 modification of the file by the asynchronous write operation.

5 In accordance with another aspect, the invention provides a file server system
6 including a first processor and a second processor for servicing client requests for access
7 to a file. The first processor has a clock for producing a clock time, and the second
8 processor has a timer for measuring a time interval. The second processor is
9 programmed for responding to a first request from a client for an asynchronous write to
10 the file by obtaining the clock time from the clock of the first processor, beginning
11 measurement of the time interval with the timer, performing a first asynchronous write
12 operation with respect to the file, and using the clock time obtained from the clock of the
13 first processor as a first file-modification time indicating a time of modification of the file
14 by the first asynchronous write operation. The second processor is programmed for
15 responding to a second request from the client for an asynchronous write to the file by
16 performing a second asynchronous write operation with respect to the file, and
17 determining a second file-modification time that is a function of the clock time obtained
18 from the clock of the first processor and the time interval measured by the timer, the
19 second file-modification time indicating a time of modification of the file by the second
20 asynchronous write operation.

21 In accordance with yet another aspect, the invention provides a file server system
22 including a first processor and a second processor for servicing client requests for access
23 to a file. The first processor has a clock for producing a clock time, and the second

1 processor has a timer for measuring a time interval. The second processor is
2 programmed for responding to a first request from a client for an asynchronous write to
3 the file by obtaining the clock time from the clock of the first processor, beginning
4 measurement of the time interval with the timer, performing a first asynchronous write
5 operation with respect to the file, and using the clock time obtained from the clock of the
6 first processor as a first file-modification time indicating a time of modification of the file
7 by the first asynchronous write operation. The second processor is further programmed
8 for receiving from the first processor an updated value for the file-modification time, for
9 comparing the updated value for the file-modification time to the first file-modification
10 time, and upon finding that the updated value for the file-modification time is greater
11 than the first file-modification time, resetting the timer. Moreover, the second processor
12 is further programmed to respond to a second request from the client for an asynchronous
13 write to the file by performing a second asynchronous write operation with respect to the
14 file, and determining a second file-modification time that is a function of the updated
15 value for the file-modification time and the time interval measured by the timer, the
16 second file-modification time indicating a time of modification of the file by the second
17 asynchronous write operation.

18 In accordance with still another aspect, the invention provides a file server system
19 including a primary processor managing metadata of a file, and a secondary processor
20 responding to requests from a client for access to the file. The primary processor has a
21 clock for producing a clock time, and the secondary processor has a timer for measuring a
22 time interval. The secondary processor is programmed for responding to a first
23 asynchronous write request from the client for writing to the file by obtaining attributes

1 of the file and the clock time from the primary processor, storing the attributes of the file
2 in a cache local to the secondary processor and using the file attributes to perform a first
3 asynchronous write operation with respect to the file, and beginning measurement of the
4 time interval with the timer. The secondary processor is further programmed for
5 responding to a second asynchronous write request from the client for writing to the file
6 by using the attributes of the file in the cache local to the secondary processor to perform
7 a second asynchronous write operation with respect to the file, and determining a file-
8 modification time that is a function of the clock time from the primary processor and the
9 time interval measured by the timer, the file-modification time indicating a time of
10 modification of the file by the second asynchronous write operation.

11 In accordance with a final aspect, the invention provides a network file server
12 including a plurality of data mover computers for servicing client requests for access to a
13 file, and a cached disk array for storing data of the file. The data mover computers are
14 coupled to the cache disk array for accessing the data of the file. The data mover
15 computers include a primary data mover computer programmed for managing metadata
16 of the file, and a secondary data mover computer programmed for requesting metadata of
17 the file from the primary data mover computer. The primary data mover computer has a
18 clock for producing a clock time, and the secondary data mover computer has a timer for
19 measuring a time interval. The secondary data mover computer is programmed for
20 responding to a first asynchronous write request from a client for writing to the file by
21 obtaining attributes of the file and the clock time from the primary data mover computer,
22 storing the attributes of the file in a cache local to the secondary data mover computer
23 and using the file attributes to perform a first asynchronous write operation with respect

1 to the file, beginning measurement of the time interval with the timer, and using the clock
2 time as a first file-modification time indicating a time of modification of the file by the
3 first asynchronous write operation. The secondary data mover computer is further
4 programmed for responding to a second asynchronous write request from the client for
5 writing to the file by using the attributes of the file in the cache local to the secondary
6 data mover computer to perform a second asynchronous write operation with respect to
7 the file, and determining a second file-modification time indicating a time of
8 modification of the file by the second asynchronous write operation.

10 **BRIEF DESCRIPTION OF THE DRAWINGS**

11 Other objects and advantages of the invention will become apparent upon reading
12 the following detailed description with reference to the accompanying drawings wherein:

13 FIG. 1 is a block diagram of a data processing system including a network file
14 server having multiple data mover computers, each of which manages a respective file
15 system;

16 FIG. 2 is a block diagram of a primary NFS server and a secondary a NFS server
17 in a file server system such as the network file server of FIG. 1;

18 FIG. 3 is a flowchart of a file attribute caching protocol between the primary NFS
19 server and the secondary NFS server of FIG. 2;

20 FIGS. 4 to 6 comprise a flowchart showing management of the file-modification
21 time attribute during the file attribute caching protocol of FIG. 3;

1 FIG. 7 is an alternative version of the flowchart of FIG. 5, showing modification
2 of the file-modification time attribute in the secondary NFS server in response to
3 notification of an update from the primary NFS server; and

4 FIG. 8 is another alternative version of the flowchart of FIG. 5, showing
5 modification of the file-modification time attribute in the secondary NFS server in
6 response to notification of an update from the primary NFS server.

7 While the invention is susceptible to various modifications and alternative forms,
8 specific embodiments thereof have been shown in the drawings and will be described in
9 detail. It should be understood, however, that it is not intended to limit the invention to
10 the particular forms shown, but on the contrary, the intention is to cover all
11 modifications, equivalents, and alternatives falling within the scope of the invention as
12 defined by the appended claims.

14 **DESCRIPTION OF ILLUSTRATIVE EMBODIMENTS**

15 In a data storage network, it is desirable to provide client access to a file system
16 through more than one processor servicing client requests. FIG. 1, for example, shows a
17 network file server that uses distributed locking and permits storage resources to be
18 incrementally added to provide sufficient storage capacity for any desired number of file
19 systems. The network file server includes multiple data mover computers 115, 116, 117,
20 each of which manages a respective file system. Each data mover computer also
21 functions as a file server for servicing client requests for access to the file systems. For
22 this purpose, each data mover computer 115, 116, 117 has a respective port to a data
23 network 111 having a number of clients including work stations 112, 113. The data

1 network 111 may include any one or more network connection technologies, such as
2 Ethernet, and communication protocols, such as TCP/IP or UDP. The work stations 112,
3 113, for example, are personal computers.

4 The preferred construction and operation of the network file server 110 is further
5 described in Vahalia et al., U.S. Patent 5,893,140 issued April 6, 1999, incorporated
6 herein by reference. The network file server 110 includes a cached disk array 114. The
7 network file server 110 is managed as a dedicated network appliance, integrated with
8 popular network operating systems in a way, which, other than its superior performance,
9 is transparent to the end user. The clustering of the data movers 115, 116, 117 as a front
10 end to the cached disk array 114 provides parallelism and scalability. Each of the data
11 movers 115, 116, 117 is a high-end commodity computer, providing the highest
12 performance appropriate for a data mover at the lowest cost. The data movers may
13 communicate with each other over a dedicated dual-redundant Ethernet connection 118.
14 The data mover computers 115, 116, and 117 may communicate with the other network
15 devices using standard file access protocols such as the Network File System (NFS) or
16 the Common Internet File System (CIFS) protocols, but the data mover computers do not
17 necessarily employ standard operating systems. For example, the network file server 110
18 is programmed with a Unix-based file system that has been adapted for rapid file access
19 and streaming of data between the cached disk array 114 and the data network 111 by any
20 one of the data mover computers 115, 116, 117.

21 In the network file server of FIG. 1, the locking information for each file system
22 119, 120, 121, 122 is managed exclusively by only one of the data movers 115, 116, 117.
23 This exclusive relationship will be referred to by saying each file system has a respective

1 data mover that is the owner of the file system. For example, the data mover 115 is the
2 owner of the A: file system 119 and the B: file system 120, the data mover 116 is the
3 owner of the C: file system 121, and the data mover 117 is the owner of D: file system
4 122. The owner of a file system is said to be primary with respect to the file system, and
5 other data movers are said to be secondary with respect to the file system.

6 In the network file server 110, each client 112, 113 may access any of the file
7 systems through any one of the data mover computers 115, 116, 117, but if the data
8 mover computer servicing the client does not own the file system to be accessed, then a
9 lock on at least a portion of the file system to be accessed must be obtained from the data
10 mover computer that owns the file system to be accessed.

11 In network file server 110, it is possible for a write operation to change the
12 attributes of a file, for example, when the extent of a file is increased by appending data
13 to the file. When a write operation will change the metadata of a file, the metadata must
14 be managed in a consistent fashion, in order to avoid conflict between the data mover
15 owning the file, and the data mover performing the write operation. For example, as
16 described in the above-cited Xu et al., U.S. Patent 6,324,581, when a secondary data
17 mover performs a write operation that changes the metadata of a file, the new metadata is
18 written to the primary data mover. This ensures that the primary data mover maintains
19 consistent metadata in its cache.

20 It is desired to permit multiple clients to have concurrent asynchronous writes to a
21 file in accordance with version 3 and version 4 of the Network File System (NFS)
22 protocol. Locking can be based on ranges of blocks within the same file. For example,
23 the primary data mover may grant one client a write lock on blocks 100 to 199 in a file,

1 and the primary data mover may grant another client a concurrent write lock on blocks
2 200 to 299 in the same file.

3 It is desirable for some of the file system metadata to be cached only on the
4 primary data mover, and some of the file system metadata to be cached on the primary
5 and secondary data movers. For example, the file system metadata is broken into three
6 categories: directory information, inodes and indirect blocks, and file attributes. For the
7 first two categories, all block allocations are performed on the primary data mover, and
8 all directory-related NFS requests are serviced on this same primary data mover.
9 However, file attributes are cached on the secondary data movers to prevent the primary
10 data mover from becoming a bottleneck for read-only access to the file attributes.

11 When multiple clients are permitted to write to the same file concurrently, it
12 becomes difficult to maintain the file-modification time attribute. Normally, when a file
13 attribute applicable to the entire file needs to be changed, the change is made at the cache
14 of the primary data mover, and the caches of the secondary data movers are invalidated.
15 The clocks of the data movers 115, 116, 177 are not synchronized. Therefore, to update
16 the file-modification time in a consistent fashion, a secondary data mover could send a
17 file-modification time request to the primary data mover, and the primary data mover
18 could read its clock to obtain a new update time, and then return the new update time to
19 the secondary data mover. Unfortunately this method would be quite burdensome,
20 because messages would have to be passed between the primary and secondary data
21 movers for each asynchronous write to a file system. In contrast, the file-creation time
22 attribute (ctime) can simply be set with the clock time of the primary data mover since a

1 file is always created by its primary data mover, and the file-creation time does not
2 change during the life of the file.

3 The file-modification time attribute must be maintained in a consistent fashion.
4 In particular, the file-modification time attribute must satisfy three important consistency
5 requirements. First, when a client writes to a file, the file-modification time should
6 increase. Second, the file-modification time should never decrease. Third, the file-
7 modification time of a file should not change unless data has actually been written to the
8 file.

9 Consistency of the file-modification time attribute is critical to the performance
10 of NFS client side caching mechanisms as well as time-based applications such as
11 incremental backup, and “make” during program compilation. If the first or second
12 consistency requirements are violated, then applications such as incremental backup and
13 “make” will become confused. If the third consistency requirement is violated, then NFS
14 clients may invalidate their cached file data unnecessarily, adversely affecting
15 performance.

16 It has been discovered that it is possible for the secondary data movers to update
17 the file-modification time attribute in a consistent fashion without always accessing the
18 primary data mover clock. The clocks of the primary and secondary data movers need
19 not be synchronized. The secondary clocks cannot simply be used to set the file-
20 modification time attribute, because the clock skew between the multiple secondary data
21 movers writing to the same file would violate the second consistency requirement. On
22 the other hand, the primary clock cannot simply be used unless the file-modification time
23 is updated for each asynchronous write. Otherwise, the third consistency requirement

1 would be violated during the gap between the time of the asynchronous write and the
2 update of the file-modification time. However, it is possible for a secondary data mover
3 to update the file-modification time attribute in a consistent fashion using a hybrid
4 method that computes the file-modification time attribute based on the clock of the
5 primary data mover and a timer of the secondary data mover. The updated file-
6 modification time is a function of the clock time obtained from the clock of the primary
7 data mover and a time interval measured by the timer of the secondary data mover.
8 Preferably, the function is a sum of the clock time obtained from the clock of the primary
9 data mover and a time interval measured by the timer of the secondary data mover.

10 As shown in FIG. 2, a first client 131 is serviced by a secondary NFS server 133,
11 and a second client 132 is serviced by primary NFS server 134. The secondary NFS
12 server 133 and the primary NFS server 134 are connected to storage 135 for access to a
13 file system 136 in the storage. For example, the NFS servers are data movers, and the
14 storage 135 is provided by a cached disk array, as described above. The secondary NFS
15 server 133 has a local cache of file attributes 137, and the primary NFS server 134 has a
16 local cache of file attributes 138. The secondary NFS server 133 has a timer 139, and the
17 primary NFS server 140 has a clock 140. The clock 140, for example, is a real-time
18 clock used by the operating system of the primary NFS server for placing a date-time
19 stamp on its local files. The timer 139 is a random access memory location that is
20 periodically incremented by a timer interrupt routine.

21 When an NFS server performs an asynchronous write for a client, the server
22 returns an updated file-modification time attribute (mtime). If the NFS server is the
23 primary NFS server 134, the updated file-modification time can simply be the time of its

1 local clock 140. If the NFS server is the secondary NFS server 133, then the updated
2 file-modification time is the sum of the local timer 139 and a local value (m) 141 of the
3 primary clock having been stored in local memory 142 of the secondary NFS server 133.
4 In particular, when a secondary 133 obtains file attributes from the primary 134 for a first
5 write to a file, the secondary receives the present value of the primary clock 140, and
6 stores the present value (m) 141 in local memory 142. At this time, the secondary resets
7 its timer 139. The secondary 133 maintains a respective timer 139 and stored clock time
8 (m) 141 for each file that it has opened for asynchronous write access.

9 When the secondary NFS server 133 performs a second asynchronous write to the
10 file system 136 for the client 131, it computes an updated file-modification time (m1) by
11 adding the stored clock time (m) 141 and the present value of its timer 139, and returns
12 the file-modification time (m1) to the client 133. When the secondary NFS server 133
13 performs a commit operation by flushing data for the file to the file system 136 in storage
14 135, the secondary NFS file server sends the updated file-modification time (m1) to the
15 primary NFS file server 134. The primary NFS file server then writes the updated file-
16 modification time (m1) to its local cache, and also sends the updated file-modification
17 time (m1) to all of the other secondaries that are caching the attributes of the file system
18 136.

19 FIG. 3 shows the preferred file management protocol (FMP) between the primary
20 NFS server and the secondary NFS server of FIG. 2. This protocol is designed to permit
21 the exchange of file metadata between primary and secondary servers that cache the file
22 metadata, as described in Xu et al., U.S. Patent 6,324,581, issued Nov. 27, 2001,
23 incorporated herein by reference. This protocol eliminates the need for the secondary to

1 communicate with the primary every time that the secondary responds to an NFS read or
2 write request from a client. In order to maintain consistency of the file attributes, the
3 primary NFS server notifies each secondary (that caches attributes for the file) whenever
4 there is a change in the attributes for the file. This allows the secondary to invalidate its
5 cache of file attributes and to refresh its cache with new attribute data from the primary.
6 In particular, in a first step 151, the secondary receives a file access request from a client.
7 In step 152, the secondary requests file attributes from the primary. The secondary does
8 this by sending a “FmpGetAttr” request 153 to the primary (over the link 118 in FIG. 1).

9 In step 154, the primary responds to the “FmpGetAttr” request from the
10 secondary by sending the file attributes 154 to the secondary and recording that the
11 secondary is caching the file attributes. In effect, the secondary is requesting a lock on a
12 range of file blocks, and if the primary can grant the range lock, then the primary returns
13 the file attributes applicable to the range of file blocks. The file attributes applicable to
14 the range of file blocks include the mapping of the logical file blocks to the logical
15 storage blocks in the storage (135 in FIG. 2). The primary may also return some file
16 attributes that apply to the entire file, such as the file’s group ID, owner, file size, file-
17 modification time (mtime) and file-creation time (ctime). In step 156 the secondary
18 receives and caches the file attributes. The secondary uses the file attributes to access the
19 file for the client. In particular, the secondary uses the mapping of the logical file blocks
20 to the logical storage blocks to read or write directly to the file system (136 in FIG. 2) in
21 the storage (135 in FIG. 2).

22 Some time later, in step 157, the primary changes the file attributes, and notifies
23 all secondaries having cached the file attributes by sending a “FmpNotify” message 158.

1 Normally, this happens only on an explicit setAttr, NFS commit, or FMP flush.
2 Therefore, an NFS asynchronous write by the client of one secondary will not result in an
3 attribute change visible to clients of another secondary. The attribute changes will be
4 visible only after a client issues an NFS commit. (This will result in the secondary
5 issuing an FMP flush.) This is consistent with NFS semantics. In step 159, the
6 secondary receives the notification, and invalidates the file attributes in its cache.

7 FIGS. 4 to 6 show management of the file-modification time attribute during the
8 file attribute caching protocol of FIG. 3. In a first step 161, the client (131 in FIG. 2)
9 initiates a first asynchronous write to a file by sending a request (WRITE3) 162 to the
10 secondary NFS server (133 in FIG. 2). In response, the secondary sends a "FmpGetAttr"
11 request 163 to the primary NFS server (134 in FIG. 2). In step 164, the primary
12 responds to the "FmpGetAttr" request by reading its clock (140 in FIG. 2) and returning
13 the file attributes and clock time (m) 165. The secondary receives the file attributes and
14 clock time (m). In step 166, the secondary stores the file attributes in its cache (137 in
15 FIG. 2) of file attributes, records the clock time (m) in its local memory (142 in FIG. 2),
16 starts its local timer (t) (139 in FIG. 2), performs the first asynchronous write to the file,
17 and caches the clock time (m) in the cache of file attributes 137 as the file-modification
18 time of the file for the first asynchronous write (WRITE3) 162. The initial value of the
19 local timer is zero. In step 167, the secondary returns file attributes including the file-
20 modification time (m) to the client.

21 Continuing in FIG. 5, in step 171, the client initiates a second asynchronous write
22 to the file. The client sends a write request (WRITE3) 172 to the secondary NFS file
23 server. In step 173, the secondary performs the second asynchronous write using file

1 attributes in its local cache (137 in FIG. 2), and calculates and caches a new value (m1)
2 for the file-modification time by adding the clock time (m) stored in its local memory
3 (142 in FIG. 2) to the value (t) of its local timer (139 in FIG. 2). This new value (m1) is
4 the file-modification time of this second asynchronous write to the file. The secondary
5 returns the file attributes and the new file-modification time (m1) 174 to the client.

6 Continuing in FIG. 6, in step 181, the client initiates a commit of the write data to
7 storage. The client sends a commit request (COMMIT3) 182 to the secondary NFS
8 server. The secondary responds by sending a flush request and the new file-modification
9 time (FmpFlush(m1)) 183 to the primary NFS server. In step 184 the primary records the
10 new file-modification time (m1) and sends it to any other clients caching attributes for
11 this file. The primary performs the requested flush operation by logging the metadata
12 changes for the client and then writing the metadata changes for the client to storage (135
13 in FIG. 2). The primary returns an acknowledgement (FlushOK) 185 to the secondary.
14 The secondary returns file attributes and the file-modification time (m1) 186 to the client.

15 It is possible for the primary to notify the secondary of a new value for the file-
16 modification time between the occurrence of the first asynchronous write and the NFS
17 commit. One way that this may happen is shown in FIG. 7, which is a modified version
18 of FIG. 5. In this case, the notification of the new value for the file-modification time
19 occurs after the second asynchronous write. Steps 191 to 194 are similar to steps 171 to
20 step 174 of FIG. 5. In step 195, the file-modification time (mp) changes on the primary
21 NFS server. The primary notifies the secondary by sending a "NotifyAttrs(mp)" message
22 196, including the new value (mp) of the file-modification time. In step 197, the
23 secondary responds by comparing the new value (mp) to its cached value (m1). If the

1 new value is not greater than its cached value, then the secondary ignores the new value
2 (mp). Otherwise, if (mp) is greater than (m1), then in step 198 the secondary caches the
3 new value (mp) as the most recent file-modification time for the file in its cache of file
4 attributes, resets its timer to zero, and sets the clock time (m) in its local memory (142 in
5 FIG. 2) to the new value (mp).

6 FIG. 8 is similar to FIG. 7 but it shows the case where the notification of the new
7 value for the file-modification time occurs before the second asynchronous write. Steps
8 201 to 204 in FIG. 8 are similar to steps 191 to 194 in FIG. 7, and steps 205 to 208 in
9 FIG. 8 are similar to steps 195 to 198 in FIG. 7.

10 Sometimes the primary might receive an FMP flush simultaneously from two
11 secondaries. In such a case, only one of the flushes will be processed. The first flush
12 processed will generate a notify message to the other client, which will invalidate the
13 server message number contained in the other client's flush. Thus the other client's flush
14 will be rejected with the error code WRONG_MSG_NUMBER.

15 The method of FIGS. 4 to 8 ensures consistency of the file-modification time
16 attribute. The first consistency requirement is met because when the client successively
17 writes to a file, the file-modification time is increased by at least the timer value (t) (in
18 step 173 of FIG. 5, step 193 in FIG. 7, and step 203 of FIG. 8).

19 The second consistency requirement is met because the sequence of file-
20 modification times on the primary server for a file is non-decreasing. In other words, if
21 m_1, m_2, \dots, m_i is the sequence of file-modification times recorded on the primary server
22 for a file, then $m_1 \leq m_2 \leq \dots \leq m_i$. This can be proven by induction on the index i .
23 For the base case of $i = 1$, the sequence is non-decreasing because it has one member m_1 .

1 For the inductive case, consider a new file-modification time m_{i+1} , which is being set on
2 the server. There are two possibilities: 1) the server received m_{i+1} from a secondary as
3 the result of an FMP flush, 2) the server received m_{i+1} locally as the result of an NFS
4 commit. For the first case 1), the secondary must have received a notification about the
5 file-modification time m_i before the flush was sent to the server (see steps 195 to 198 in
6 FIG. 7). At that time, the secondary compared m_i to its current in-memory file-
7 modification time $m1 = m_x + t$ (see step 197 in FIG. 7 and step 207 in FIG. 8). If m_i was
8 greater, the secondary used m_i as the new basis $m1$ for its file-modification time and reset
9 its timer (step 198 of FIG. 7 and step 208 of FIG. 8); otherwise it ignored m_i . Let d be
10 the delta between the receipt of m_i at the secondary, and the last asynchronous write at
11 the secondary before the fmp flush. Then

12 $m_i > m_x + t$ implies $m_{i+1} = m_i + d$, and

13 $m_i \leq m_x + t$ implies $m_{i+1} = m_x + t + d$.

14 Because d is greater than or equal to zero, we conclude $m_{i+1} \geq m_i$.

15 For the second case 2), the argument is the same, because when the
16 primary notifies other secondaries of a new file-modification time m_i for a file, it
17 also checks its own local in-memory file-modification time $m1$, and if the local
18 time is behind m_i , then its in-memory file-modification time is set to m_i .

19 The third consistency requirement is met because the method of FIGS. 4-8 only
20 changes the file-modification time of a file unless data has actually been written to the
21 file. For example, if a client issues an NFS write, and receives m_i as the file-
22 modification time in the post-op attributes, and then issues an NFS commit, it will be

1 guaranteed to see m_i as the file-modification time in the post-op attributes, unless another
2 secondary has issued an FMP flush to the server in the meantime. If such a flush was
3 issued, then the NFS client should invalidate its cache, so this behavior is not
4 problematic. The key fact is that the act of issuing the NFS commit to the secondary
5 server in and of itself does not change the file-modification time of the file, from the
6 point of view of the NFS client. This ensures that the management of the file-
7 modification time will not cause problems for NFS client caching schemes.

8 It should be apparent that the structure and operation shown in FIGS. 1 to 8 can
9 be modified in various ways that are covered by the appended claims. For example, the
10 NFS servers 133, 134 in FIG. 2 could be geographically remote from each other and
11 remote from the storage 135 and interconnected in a wide-area data network. In addition,
12 the timer 139 in the secondary NFS server 133 could be reset with the clock time (m)
13 from the primary NFS server 134 (or with the updated value (mp) for the file-
14 modification time) instead of being reset to zero, so that the timer 139 would periodically
15 compute a sum of the clock time (m) (or the updated value (mp)) and the time interval (t)
16 measured by the timer.

17 In view of the above, there has been described a method of maintaining a file-
18 modified time attribute in a multi-processor file server system. To permit multiple
19 unsynchronized processors to update the file-modification time attribute of a file during
20 concurrent asynchronous writes to the file, a primary processor manages access to
21 metadata of the file, and has a clock producing a clock time. A number of secondary
22 processors service client request for access to the file. Each secondary processor has a
23 timer. When the primary processor grants a range lock upon the file to a secondary, it

1 returns its clock time (m). Upon receipt, the secondary starts a local timer (t). When the
2 secondary modifies the file data, it determines a file-modification time that is a function
3 of the clock time and the timer interval, such as a sum (m+t). When the secondary
4 receives an updated file-modification time (mp) from the primary, if $mp > m+t$, then the
5 secondary updates the clock time (m) to (mp) and resets its local timer.

6 Although the method of maintaining the file-modified time attribute has been
7 described above with respect to a network file server as shown in FIG. 1 or FIG. 2, it
8 should be understood that the method has general applicability to diverse kinds of file
9 server systems, such as server clusters and storage area networks. For example, when it
10 is desired to permit more than one processor in such a system to change the file-modified
11 time attribute of a file, the method can be used to eliminate a need to synchronize the
12 processors or to require the processors to always obtain the file-modified time attribute
13 from a common clock.

14

15